

```

1  --
2  --
3  --
4  --
5  --
6  --
7  --
8  --
9  local l = require"lib"
10 local the = l.settings[[
11
12 RL.LUA : stings
13 (c)2022 Tim Menzies <tim@ieee.org> BSD(2clause).
14
15 USAGE:
16   lua rlgo.lua -[bghk] [ARG]
17
18 OPTIONS:
19   -b --bins      discretization control = 8
20   -F --Far       in "far", how far to seek = .95
21   -g --go        start-up action      = pass
22   -h --help      show help
23   -k --keep      keep only these nums = 256
24   -s --seed      random number seed   = 10019
25   -S --Some      in "far", how many to search = 512]]
26
27 local About= {} -- factor for making columns
28 local Col = {} -- summarize one column
29 local Data = {} -- store rows, and their column summaries
30 local Row = {} -- store one row
31
32 -- CODE CONVENTIONS
33 -- Leading_upper_case : class
34 -- l.      : instance va
35 -- l.s     : reference to a library function
36 -- prefix _ : some internal function,variable.
37
38 -- type hints: where practical, on function arguments,
39 --   - t = table
40 --   - prefix s=string
41 --   - prefix n=num
42 --   - prefix is=boolean
43 --   - class names in lower case denote vars of that class
44 --   - suffix s denotes table of things
45
46 -----
47 -- About
48 --
49 --
50 function About.new(sNames)
51   return About._cols{[names=sNames, all={}, x={}, y={}, klass=nil},sNames} end
52
53 -- How to recognize different column types
54 local _is={
55   nom   = "[a-z]", -- ratio cols start with uppercase
56   goal  = "[!a-z]", -- !=klass, {r,-}=maximize,minimize
57   klass = "[s]", -- klass if "+"
58   skip  = "[*]", -- skip if "*"
59   less  = "-$]" -- minimize if "-"
60
61 -- Turn a list of column names into Col objects. If the new col is independent
62 -- or dependent of a goal attribute then remember that in i.x or i.y or i.klass.
63 function About._cols(i,sNames)
64   for at,name in pairs(sNames) do
65     local col = l.push(i.all, Col.new(name,at))
66     if not name:find(_is.skip) then
67       l.push(name:find(_is.goal) and i.y or i.x, col)
68       if name:find(_is.klass) then i.klass=col end end end
69   return i end
70
71 -- Update, only the non-skipped cols (i.e. those found in i.x and j.x.
72 function About.add(i,t)
73   local row = t.cells and t or Row.new(i.about, t)
74   for _,cols in pairs(i.x,i.y) do
75     for _,coll in pairs(cols) do
76       Col.add(coll, row.cells[coll.at]) end end
77   return row end
78
79 -----
80 -- Col
81 --
82 --
83 -- Summarize one column.
84 function Col.new(txt,at)
85   txt = txt or ""
86   return {n = 0, -- how many items seen?
87     at = at or 0, -- position of column
88     txt = txt, -- column header
89     isNom = txt:find(_is.nom), -- and -1 or 1,
90     v = txt:find(_is.less) -- false if some update needed
91     ok = true, -- false if some update needed
92     _has = {} end -- place to keep (some) column values.
93
94 -- Update
95 function Col.add(i,x)
96   if x ~= "?" then
97     i.n = i.n + 1
98     if i.isNom
99       then i._has[x] = 1 + (i._has[x] or 0)
100     else local pos
101       if #i._has < the.keep then pos= 1 + (#i._has)
102       elseif l.rand() < the.keep/i.n then pos=l.rand(#i._has) end
103       if pos then
104         i.ok=false -- kept items are no longer sorted
105         i._has[pos]=x end end end end
106
107 -- Distance
108 function Col.dist(i,x,y)
109   if x=="?" and y=="?" then return l end
110   if i.isNom
111     then return x==y and 0 or 1
112   else if x=="?" and y=="?" then return l end
113   if x=="?" then y = Col.norm(i,y); x=y<.5 and 1 or 0
114   elseif y=="?" then x = Col.norm(i,x); y=x<.5 and 1 or 0
115   else x,y = Col.norm(i,x), Col.norm(i,y) end
116   return math.abs(x-y) end end
117
118 -- Diversity
119 function Col.div(i)
120   if i.isNom
121     then local e=0
122     for _,v in pairs(i._has) do
123       if v>0 then e=e-v/i.n*math.log(v/i.n,2) end end

```

```

126   return e
127   else local t=Col.has(i)
128   return (l.per(t,.9) - l.per(t,.1))/2.56 end end
129
130 -- Sorted contents
131 function Col.has(i)
132   if i.isNom then return i._has end
133   if not i.ok then table.sort(i._has) end
134   i.ok=true
135   return i._has end
136
137 -- Central tendency
138 function Col.mid(i)
139   if i.isNom
140     then local mode,most=nil,1
141     for k,v in pairs(i._has) do if v>most then mode,most=k,v end end
142     return mode
143   else return l.per(Col.has(i),.5) end end
144
145 -- Return num, scaled to 0..1 for lo..hi
146 function Col.norm(i,num)
147   local a= Col.has(i) -- "a" contains all our numbers, sorted.
148   return a[#a] - a[1] < 1E-9 and 0 or (num-a[1])/(a[#a]-a[1]) end
149
150 -- Map x to a small range of values.
151 function Col.discretize(i,x, a,b,lo,hi)
152   if i.isNom then return x else
153     a = has(i)
154     lo,hi = a[1], a[#a]
155     b = (hi - lo)/the.bins
156     return hi==lo and 1 or math.floor(x/b+.5)*b end end
157
158 -----
159 -- Row
160 --
161 --
162 -- Hold one record
163 function Row.new(about,t)
164   return {_about=about, cells=t, cooked=l.map(t,l.same)} end
165
166 -- Everything in rows, sorted by distance to i.
167 function Row.around(i,rows)
168   local fun = function(j) return (row=j, d=Row.dist(i,j)) end
169   return l.sort(l.map(rows, fun), lt"d") end
170
171 -- Recommend sorting i before j (since i is better).
172 function Row.better(i,j)
173   i.evald,j.evald= true,true
174   local s1,s2,d,n,x,y=0,0,0,0
175   local ys,e = i._about.y,math.exp(1)
176   for _,col in pairs(ys) do
177     x,y= i.cells[col.at], j.cells[col.at]
178     x,y= Col.norm(col,x), Col.norm(col,y)
179     s1 = s1 - e^(Col.w * (x-y)/ys)
180     s2 = s2 - e^(Col.w * (y-y)/ys) end
181   return s1/ys < s2/ys end
182
183 -- Distance
184 function Row.dist(i,j)
185   local d,n,x,y,distl=0,0
186   local cols = cols or i._about.x
187   for _,col in pairs(cols) do
188     x,y = i.cells[col.at], j.cells[col.at]
189     d = d + Col.dist(col,x,y)^the.p
190     n = n + 1 end
191   return (d/n)^(1/the.p) end
192
193 -----
194 -- Data
195 --
196 --
197 -- Holds n records
198 function Data.new(t) return {rows={}, about=About.new(t) } end
199
200 -- Update
201 function Data.add(i,t) l.push(i.rows, About.add(i.about,t)) end
202
203 -- Replicate structure
204 function Data.clone(i, t)
205   local out = Data.new(i.about.names)
206   for _,row in pairs(t or {}) do Data.add(data,row) end
207   return data end
208
209 -- Discretize all row values (writing those vals to "cooked").
210 function Data.discretize(i)
211   for _,col in pairs(i.about.x) do
212     for _,row in pairs(i.rows) do
213       local x = row.cells[col.at]
214       if x ~= "?" then
215         row.cooked[col.at] = discretize(col,x) end end end end
216
217 -- Recursively bi-cluster one Data into sub-Datas.
218 function Data.cluster(i, rowAbove,stop)
219   stop = stop or (#i.rows)^the.Min
220   if #i.rows >= 2*stop then
221     local A,B,As,Bs,c = Data.half(i.rows,rowAbove)
222     i.halves = {c=c, A=A, B=B,
223       kids = { Data.cluster(Data.clone(i,As), A, stop),
224         Data.cluster(Data.clone(i,Bs), B, stop) }}end
225   return i end
226
227 -- Split data according to distance to two distant points A,B
228 -- To speed things up, find distant points via A=Far(any()) and B=Far(A).
229 -- To speed things up, try to reuse a distant point from above (see RowAbove).
230 -- To speed things up, only look at some of the rows (see the.See).
231 -- To dodge outliers, don't search all the way to edge (see the.Far).
232 function Data.half(i, rows, rowAbove)
233   local some= l.many(rows, the.See)
234   local function far(row)
235     return l.per(Row.around(row,some), the.Far).row end
236   local function project(row)
237     local a,b = Row.dist(row,A), Row.dist(row,B)
238     return (row=row, x=(a^2 + b^2)/(2*c)) end
239   local A= rowAbove or far(l.any(some))
240   local B= far(A)
241   local c= Row.dist(A,B)
242   local As,Bs = {}/,{}
243   for n,rowx in pairs(l.sort(l.map(rows, project),l.lt"*")) do
244     push(n < #rows/2 and As or Bs, rowx.row) end
245   return A,B,As,Bs,c end
246
247 -- Load from file
248 function Data.load(sFilename, data)

```

```

251   l.csv(sFilename, function(row)
252     if data then Data.add(data,row) else data=Data.new(row) end end)
253   return data end
254
255 -- Central tendency
256 function Data.mid(i) return l.map(i.about.y, Col.mid) end
257
258 -- Guess the sort order of the rows by peeking at a few distant points.
259 function Data.optimize(i, rowAbove,stop,out)
260   stop = stop or (#i.rows)^the.Min
261   out = out or {}
262   if #i.rows < 2*stop
263     then for _,row in pairs(i.rows) do push(out,row) end
264     else local A,B,As,Bs,c = Data.half(i.rows, rowAbove)
265     if Row.better(A,B)
266       then for j=#Bs,1 do push(out,Bs[j]) end
267       Data.optimize(Data.clone(i,rev(As)), A, stop, out)
268     else for _,row in pairs(As) do push(out,row) end
269     Data.optimize(Data.clone(i,Bs), B, stop, out)
270   end end
271   return out end
272
273 -----
274 return {Data=Data,Row=Row,Col=Col,About=About,the=the}

```